

(19) World Intellectual Property Organization  
International Bureau(43) International Publication Date  
1 March 2001 (01.03.2001)

PCT

(10) International Publication Number  
WO 01/14961 A2(51) International Patent Classification<sup>7</sup>: G06F 9/00

(21) International Application Number: PCT/US00/24336

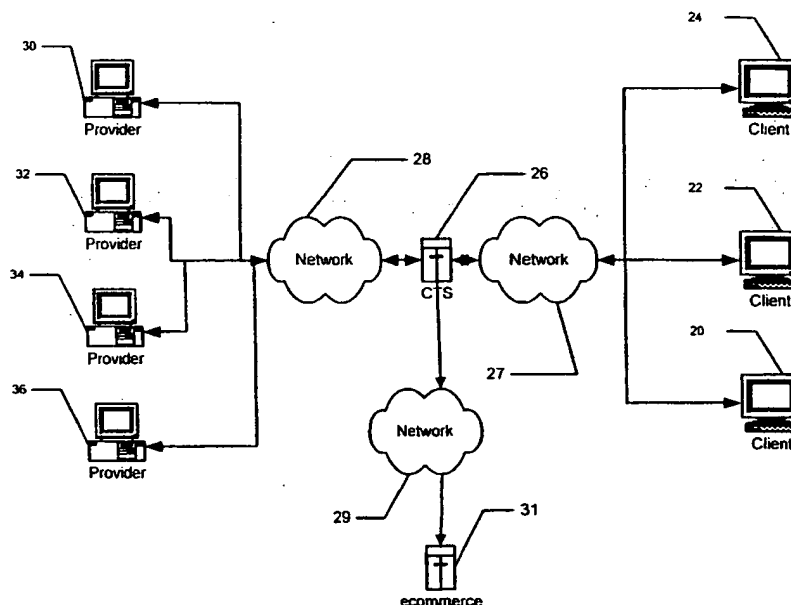
(22) International Filing Date: 28 August 2000 (28.08.2000)

(25) Filing Language: English

(26) Publication Language: English

(30) Priority Data:  
60/150,766 26 August 1999 (26.08.1999) US  
60/210,334 12 June 2000 (12.06.2000) US(71) Applicant: PARABON COMPUTATION [US/US];  
Suite 1000, 3930 Walnut Street, Fairfax, VA 22030 (US).(72) Inventors: ARMENTROUT, Steven, L.; 11879 St.  
Trinians Ct., Reston, VA 20191 (US). O'CONNOR,  
James; 11625 Ayreshire Road, Oakton, VA 22124 (US).  
GANNON, James; 3817 Inverness Road, Fairfax, VA  
22030 (US). SLETTEN, Brian; Apt. 906, 900 N. Stuart  
Street, Arlington, VA 22203 (US). CIER, Sean; Apt.12, 11800 Federalist Way, Fairfax, VA 22030 (US).  
CARLSON, Sarah; 13236 Memory Lane, Fairfax, VA  
22033 (US). DAVIS, Jonathan; Apt. 312, 202 Columbia  
Road NW, Washington, DC 20009 (US). DUPERTUIS,  
Greg; 7888 Leroux Lane, Manassas, VA 20112 (US).  
MCCLOUGHLIN, Scott; 2827 28th Street, #14, Wash-  
ington, DC 20001 (US). DAVIES, Antony; 11733 Stuart  
Mill Road, Oakton, VA 22124 (US).(74) Agents: ROBERTS, Jon, L. et al.; Roberts Abokhair &  
Mardula, LLC, Suite 1000, 11800 Sunrise Valley Drive,  
Reston, VA 20191 (US).(81) Designated States (*national*): AE, AG, AL, AM, AT, AU,  
AZ, BA, BB, BG, BR, BY, BZ, CA, CH, CN, CR, CU, CZ,  
DE, DK, DM, DZ, EE, ES, FI, GB, GD, GE, GH, GM, HR,  
HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR,  
LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, MZ,  
NO, NZ, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM,  
TR, TT, TZ, UA, UG, UZ, VN, YU, ZA, ZW.(84) Designated States (*regional*): ARIPO patent (GH, GM,  
KE, LS, MW, MZ, SD, SL, SZ, TZ, UG, ZW), Eurasian  
patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European

[Continued on next page]

(54) Title: SYSTEM AND METHOD FOR THE ESTABLISHMENT AND UTILIZATION OF NETWORKED IDLE COMPU-  
TATIONAL PROCESSING POWER

(57) Abstract: A distributed computing platform using the idle computational processing power of a plurality of provider computers is disclosed. At least one networked server collects tasks from client computers, schedules and distributes the tasks to networked provider computers, and collects and returns results to client computers. A client API forms tasks and collects results. A compute engine operates on the provider computers to communicate with the server and execute tasks using idle computational power.

WO 01/14961 A2

1   **Title of Invention:**   System and Method for the Establishment and Utilization of Networked  
2                                   Idle Computational Processing Power  
3

4   **Field of the Invention:**

5               This invention relates generally to computer use over a network to form a distributed  
6   computing platform. More specifically the present invention is a system and method for use  
7   of idle time of computers that are connected to a network by users requiring significant  
8   computing power without the need for a large scale dedicated processing unit.  
9

10   **Background of the Invention:**

11              When a computer is operating, but not actively performing computations for someone,  
12   it is said to be *idle*. Because of their incredible speed, modern computers are idle most of the  
13   time, not only when they are running screen savers, but *even when they are being used*. For  
14   instance, a fast typist working at top speed in a word processor uses only a fraction of the  
15   available computational capacity of a desktop PC. Although the time between a typist's  
16   keystrokes seems immeasurably small to humans, it is an eternity when measured in computer  
17   time - an eternity filled with unproductive *idle computation*. Because of this gross  
18   underutilization, it is estimated that well over 95% of the world's computational capacity  
19   presently goes to waste.  
20

21   **Summary of the Invention**

22              The present invention is drawn to organization and use of the idle processing power of  
23   general-purpose computers to form a distributed computing platform.

24              It is therefore an objective of the present invention to provide large amounts of  
25   computational power to users without the users having to purchase a large computer for such  
26   purposes.

27              It is a further objective of the present invention to harness the idle computational  
28   power of many computers and make that power available to clients on an as needed basis.

1           It is yet another objective of the present invention to allow clients to create a virtual  
2 cluster of machines of client-definable computational power to run client tasks.

3           It is thus an overall goal of the present invention to utilize this heretofore underutilized  
4 computer power in a novel way as a distributed computing platform to meet the needs of users  
5 who require vast computing power but who may not have the financial wherewithal to  
6 purchase or lease large mainframe computers or other supercomputing solutions.

7           Just as email service providers deliver email from one user to another, a business  
8 running the system and method of the present invention will disseminate the tasks of a client's  
9 distributed computer program to multiple providers for remote execution. As the remote  
10 providers complete tasks, the providers will deliver results back to the originating client. As  
11 part of the present invention, an intermediary server delivers *computation* from providers to  
12 clients.

13           The Internet infrastructure to facilitate the distributed processing of the present  
14 invention requires three components:

- 15           1) Client software applications for requesting and submitting distributed processing  
16 requests and jobs. This is accomplished through a client API that allows client  
17 jobs to be coded as independent tasks. The underlying task model, while complex,  
18 is hidden from the user.;
- 19           2) Provider software called a compute engine (CE) for managing the launch and  
20 execution of tasks delivered from the intermediary server. The CE runs  
21 unobtrusively and securely on the provider's computer. It processes tasks when  
22 the computer is idle and returns results to the server when the provider computer  
23 connects to the network which, for example, and without limitation may be the  
24 Internet; and
- 25           3) A centralized task server for exchanging tasks and results between participants.  
26 The centralized task server (CTS) comprises a three-tiered architecture as more  
27 fully set forth below and is redundant, fault tolerant and scalable as traffic and  
28 clients increase. A task scheduler in the CTS matches power requested to  
29 available provider computer resources.

1 In this fashion, from a client's perspective, the present invention represents a gigantic virtual  
2 computer or distributed computing platform, ideally suited for performing large-scale parallel  
3 *distributed computation*, a function that was formerly the exclusive domain of traditional  
4 "chips-and-solder" supercomputers. The present invention, however, will perform such  
5 computations at a fraction of the cost and on an as-needed basis.

6 The business model associated with the present invention is to create a new market for  
7 idle computation power. A business running the apparatus and method of the present  
8 invention will purchase the idle computational power of millions of *providers* (individual  
9 computer owners), and deliver that idle computational power as a single entity to *clients*  
10 (technologists seeking ultra-high-performance computation).

11 Presently, idle computation has no value whatsoever because the market for  
12 computation is inseparable from the market for computers. Users have no means of buying,  
13 for example, larger amounts of times on disparate computers to execute a job when the user  
14 needs it and not just when a specific computer is available.

15 With the present invention, however, immense value can be created by a  
16 computational intermediary in the same way that financial engineers create value by  
17 constructing derivative financial instruments. In effect, the present invention *strips* excess  
18 computational capacity from provider computers, *bundles it* into quantities suitable for large-  
19 scale parallel computations, and *delivers* it to clients on an as needed basis. Providers benefit  
20 by selling a heretofore-wasted resource. Clients benefit because they can inexpensively  
21 purchase just the quantity of computation needed, thereby avoiding the normally huge fixed  
22 costs associated with purchasing a supercomputer.

23 The business model of the present invention creates and intermediates the new market  
24 in idle computation.

25 Finally, an applications research and development group will add value by utilizing its  
26 superior knowledge of distributed processing to solve high-payoff computational problems.  
27 Such a team will spur demand in two ways: First, its purchases of idle computation will  
28 stimulate market activity; second, the success of applications and research group will  
29 demonstrate to prospective clients the power and cost-effectiveness of idle time processing.

1       A centralized task server (CTS) 26 intermediates the needs of client applications and  
2 the power of the provider computers 30-36 to meet the needs of clients for computational  
3 power. Clients make requests via client computers 20-24 and by virtue of a client application  
4 programming interface (Client API) that allows the client to specify the amount of computing  
5 power desired, the individual rating of computers that will execute the client's job (e.g., 2  
6 gigaflops or GF [Flops = Floating point operations per second]), how much memory and disk  
7 space will be necessary and other parameters essential to running the client's job. The CTS  
8 receives tasks from clients and assigns them to one or more providers based the characteristics  
9 of the available providers and the job characteristics specified by the client. When providers  
10 return results, the CTS forwards these back to the client application.

11       Each provider computer has a resident compute engine (CE) that receives and  
12 executes the various tasks that are assigned to it by the CTS. The CE runs client tasks and  
13 manages the use of idle computation power of the provider computer. The CE returns  
14 intermediate and final results to the CTS that then forwards them back to the client  
15 application.

16       Communication between the client application, CTS, and CE is governed by a  
17 communication protocol unique to the present invention. At a high level of abstraction,  
18 separate message and data channels exist. The message channel is a bi-directional channel  
19 between clients and the server that is used to pass high-level messages that facilitate  
20 distributed processing. The data channel is a channel used to pass large blocks of data called  
21 data elements and executable elements. Executable elements are blocks code developed by  
22 clients or a third-party. This code represents tasks or portions of tasks. Data elements are  
23 blocks of data that are either input to or output of client tasks. High-level messages  
24 transmitted on the message channel contain references to data and executable elements.  
25 These items are separately transferred on the data channel on an as needed basis.

26       The details of the client development environment, CTS, CE, and protocol the  
27 protocol that these components use to communicate will be covered in subsequent sections.

28       Referring first to **Figure 2** the general flow of the present invention is illustrated.  
29 Clients 10, that is, those entities desiring execution of computationally intensive tasks, launch

1 allows a client application to manage the execution of jobs consisting of a very large number  
2 of tasks, and a set of tools that support the monitoring and control of very large jobs. These  
3 are described in more detail below.

4 The runtime API defines the environment in which client task run. The runtime API  
5 consists of a set of function calls that the client task will exercise during execution to interact  
6 with its environment. This API contains calls that the task can use to:

- 7 • Obtain access to the initialization parameters for the task.
- 8 • Obtain access to data elements needed to perform its function.
- 9 • Perform a checkpoint so that if the task is interrupted, it can restart at the checkpoint  
10 and therefore reduce the loss of work.
- 11 • Send interim and final results back to the server.

12 In addition, the API defines entry points that allow the environment to control the task.

13 These include:

- 14 • Start a task (note that restarting a task is a special case of starting a task).
- 15 • Stop (suspend) a task.
- 16 • Request a checkpoint.
- 17 • Request a status.
- 18 • Terminate a task.

19 This is not an exhaustive list of functions and is only illustrative of what can be  
20 accomplished via the runtime API.

21 There can be multiple implementations of this API which allow the task to run in  
22 multiple environments. The two most important implementations for the purpose of this  
23 invention are the implementation in the CE and the run local implementation included in the  
24 client development environment. The CE implementation allows the task to run in the CE on  
25 a provider's box (regardless of platform). The run local implementation allows clients to  
26 execute tasks on their local client machine. This enables debugging and testing before the job  
27 is distributed. However, the client could elect to run an entire job (or a certain number or  
28 percent of tasks) locally.

1 (percentage), the work done on completed tasks, and the time spent per completed task.  
2 These statistics are illustrative only. Other statistics may also be gathered as desired by the  
3 client. The information that is monitored is displayed in both graphic and text form.

4 These tools can also be used to support other general-purpose functions such as  
5 starting and stopping individual tasks or jobs, removing individual tasks or jobs, etc.

6 Referring to **Figure 3**, one embodiment of the CE architecture of the present invention  
7 is illustrated. The CE is the software unit that resides on a provider computer and  
8 accomplishes the timely and secure execution of client tasks. The CE 84 manages the  
9 interface between the provider computer and the server of the present invention, monitors the  
10 state of the provider computer's machine, schedules the execution of tasks on the provider  
11 computer, and manages all communication with the server of the present invention. For  
12 example, the CE 84 receives tasks from the CTS, posts results to the server, provides status of  
13 tasks and receives messages to remove, suspend or resume tasks.

14 The CE 84 can be configured to execute tasks continuously, albeit at low priority (so  
15 as not to interfere with a provider's normal use of their computer), or to start execution of  
16 tasks when a provider's computer is idle (i.e., not being interactively used) and stop execution  
17 immediately upon provider interaction with their computer. In this way, the CE is designed to  
18 be unobtrusive to the provider.

19 CE 84 communications with the server and execution of tasks occur independently.  
20 For example, even if a provider is actively using their computer, the CE communicates with  
21 the server, provided a network connection is available (e.g., when a provider has dialed into  
22 the Internet). In so doing, the CE receives elements that are placed in the task queue 86 to be  
23 executed whenever possible. Thus, obtaining and staging of client tasks takes place at a time  
24 independent from the execution of those tasks. Making CE communications independent  
25 from CE task execution adds to the efficiency of the overall system.

26 The CE also includes a dialing scheduler that can be preset by the provider to dial out  
27 to the server of the present invention at configurable intervals and only during preset periods.  
28 Thus, for example, the provider can direct the dialer to dial out to the server every ten  
29 minutes, but only between the hours of 2 a.m. to 4 a.m.

1 accessible to the user such that accounts can be accessed, amount of time associated with use  
2 of the users computer can be determined.

3 The graphic window of the CE is displayed for the provider when the CE detects the  
4 host computer is idle. Its presence indicates to the provider that the system is processing, or is  
5 ready to process tasks in the queue, and has a number of functions. The window can show the  
6 state of processing, that is whether tasks are being processed at a particular moment or not.  
7 Further, within the window are located various other buttons that may be clicked to determine  
8 the status of the provider's account. In addition, advertising or other messages may be  
9 presented in the CE window, which will allow a user to proceed to a particular website to  
10 view further information concerning the messages. Thus, the CE window can operate in the  
11 nature of a browser.

12 The size and placement of the CE window during "pop up" is configurable by  
13 the user. When the CE window is displayed, the system seizes control of the cursor and  
14 locates it at any configurable spot within the CE window. Thus the cursor can be placed in a  
15 spot where messages are to be displayed or in a location to actuate any functionality that is  
16 available to the provider while the CE is operating, such as the quick disappearance of the  
17 display window as described above.

18 The user can also confine the CE to a system tray if desired. In this case, the icon or  
19 panel in the tray can be used to indicate, among other things, engine status, the arrival of new  
20 display content or upgrade packages, or Internet connectivity status. Alternatively, the CE  
21 can be configured and executed as a "daemon" having only a command-line user interface.

22 The CE also supports the dynamic display of marketing content in the CE window to  
23 the provider computer. The server notifies the CE when new content is available and where it  
24 can be downloaded.

25 In the event that the user desires to use the computer for purposes other than to run  
26 elements of the present invention, the user will simply type a key or cause the cursor to be  
27 moved outside the CE window. As soon as this occurs, the CE will shut down all of the  
28 processing of the present invention in an orderly fashion. Alternatively, the CE can be  
29 configured to execute tasks constantly, albeit at low priority, in which case the CE window



1           The runtime then requests the data elements for the task **106** over the runtime to CE  
2 channel, and the CE returns them in response over the CE to runtime channel.

3           The runtime then establishes a task context that the task uses to communicate with the  
4 runtime system **108**.

5           The runtime establishes and begins the client task **110**. The executable elements for  
6 this task having been provided to the runtime via the class path provided for the JVM when  
7 the JVM was started. The client task then runs in a separate protection domain that prevents it  
8 from accessing the network or the disk. This protection domain is established through the  
9 permissions granted to the task in the task policy file.

10           The client task interacts with the runtime through the task context object established  
11 for it by the runtime **112**. The runtime calls run with the client tasks privileges when possible.  
12 If the runtime needs to perform an operation requiring additional privileges, i.e., accessing a  
13 particular data element, the runtime executes a privilege block that takes on the minimum  
14 privilege for the minimum amount of time required to complete the operation.

15           The client interacts with the runtime to get task parameters and to access task elements  
16 **114**. As an optional step, when the client task can send out temporary results using the set  
17 status call to the runtime **116**, these results are passed to the CE over the runtime to CE  
18 channel. The CE then sends these interim results in the form of a task status to the server.

19           As an additional option, when the client task can create a task checkpoint using the  
20 checkpoint call to the runtime **118**, the checkpoint is passed to the CE over the runtime to CE  
21 channel. The CE serializes the checkpoint to disk so it is available in case a task restart is  
22 required.

23           When finished, the client task does a set status runtime call **120** with the final results  
24 and exits. The results are passed through the CE over the runtime to CE channel. The CE  
25 then sends these results in the form of a task status message to the server. After the CE sends  
26 a final task status noting that the task is complete **120**, the space that is consumed by the task  
27 definition and related data can be reclaimed for use in running other tasks.

28           Once the task completes **122**, the runtime closes down its connections and exits.

1           One of the most critical functions of the CTS is task scheduling. Part of the  
2 scheduling process is to provide clients of the present invention, who desire to obtain the  
3 aggregate computational power, with the desired amount of processing power on demand.

4           The server of the present invention loads tasks onto a single provider computer based  
5 in such a way as to match computer capability with task requirements. As tasks are  
6 completed, and through the communications channel between the provider computer and the  
7 server, additional tasks are fed to the provider computer for execution. All new tasks are  
8 placed in a task execution queue on the provider's computer for subsequent execution. To  
9 provide maximal control of the execution of tasks on a provider's computer, the CE supports  
10 the concept of task priorities and task preemption (tasks can preempt executing tasks of lower  
11 priority). These features allow the server to more effectively utilize available provider  
12 computer resources.

13           In order for the server of the present invention to correctly characterize the provider's  
14 computer, the server also collects usage and performance histories (profiles) from the provider  
15 computer concerning a wide variety of parameters. For example, the server collects Internet  
16 connectivity patterns and processor availability patterns during the course of any 24-hour  
17 period, weekly period, monthly period, etc. The processing power of the processor itself is  
18 also important since certain clients may request higher or lower processing power depending  
19 on relative cost and timeliness tradeoffs.

20           To fulfill its function, the task scheduler must fulfill the following requirements: The  
21 task scheduler gets jobs and associated tasks from clients. It must then take those tasks and  
22 distribute them to various provider computers for execution. It must also record the fact that  
23 it has assigned a particular task to a particular provider computer.

24           The task scheduler passes both interim and final results to the client and records the  
25 amount of work performed to an on-line ledger. It then credits a provider's account with the  
26 amount of time that the provider computer actually performed on the specific task. This  
27 allows the provider to be paid by the computational work that it has contributed to a particular  
28 job. (By definition, time multiplied by the power of the provider computer equals work  
29 performed.)

1 provider computer and the bandwidth that is available to provide elements and tasks to  
2 provider computers.

3 The CTS collects statistics regarding provider computer capabilities. This is critical  
4 for both assigning payment to a provider and a debit notice to the client. It is also important  
5 that the CTS know at all times the capability of the provider computers.

6 As noted above, the CTS may send the same tasks to more than one provider  
7 computer. This helps the CTS be assured that the computations that are performed by the  
8 assigned provider computers are correct. Since results from provider computers should be the  
9 same for the same task (unless randomization is integral to the calculation of the task), any  
10 differences in response are noted by the CTS and resolved. This assures the client that  
11 answers are accurate.

12 The CTS takes explicit measures to protect the intellectual property inherent in a task.  
13 To accomplish this, the CTS obfuscates executable elements and data elements to remove  
14 contextual clues to the code's intended purpose. Furthermore, the CTS also removes from  
15 tasks all traces of client identity. Thus, unless specifically intended otherwise by the CTS, a  
16 provider has no means of knowing the clients for whom elements are executed nor their  
17 domain of inquiry.

18 To characterize and assess the capabilities of provider computers, in order to perform  
19 task scheduling effectively, the CTS executes certain benchmark tasks on provider computers.  
20 Also, various test tasks with known results are sent to the provider computer to ensure that  
21 results returned by the provider computer are accurate. Unexpected results from a provider  
22 can be used to identify attempts to submit fraudulent results or otherwise tamper with the  
23 normal execution of the system of the present invention.

24 In operation, the present invention allows the client to specify job and task parameters  
25 via the client interface. Using the client interface, the Client selects desired attributes of  
26 nodes in a virtual cluster. These attributes include, but are not limited to, CPU speed (in GF),  
27 memory, available disk space, intermittent or constant connectivity, bandwidth (in Kbps),  
28 number of nodes, and reliability (as measured by consistency of results). The client also

$$\text{Price per hour for computation} = 4.5 + 0.296P^{1.155}$$

where 4.5, 0.296 and 1.155 are configurable constants.

The client is charged an additional amount per GB for data transferred through each of the four communications paths, 150, 152, 154, and 156 of the present invention as illustrated in Figure 7. The communications paths 150, 152, and 154 represent duplicated data and are a function of the degree of reliability the client achieves. Of the dashed arrows, only one will be executed, and so the client will be charged for only one of these paths.

The client is also charged for data transfer during the course of executing the particular tasks of the client. For example, let  $d_i$  be the quantity (in GB) of data transferred from the client to the server. Let  $d_o$  be the quantity of data transferred from the server to the client. Given that some of the data will be duplicated for processing by virtual nodes comprised of more than one provider computer, the client's total data transfer charge is given by the following:

$$\text{Total Cost of Data Transfer} = g(d_i + d_o) \left( \frac{\ln(1-r)}{\ln(1-\pi)} + 1 \right)$$

where  $g$  is the transfer cost per leg per GB,  $r$  is the requested node reliability, and  $\pi$  is the average provider computer reliability.

Requested node reliability  $r$  can be defined for a virtual node of CPU speed  $C$  and bandwidth  $B$  as the probability of the virtual node completing a task in no more time than the time it would take a fully idle, constantly connected computer with CPU speed  $C$  and bandwidth  $B$  to complete the task.

Average provider computer reliability  $\pi$  is defined as the probability of the provider's frequency of contact not diminishing, in the near term, from the provider's historic average.

The client is also charged an hourly premium for requesting reliability greater than a set baseline reliability. The charge per hour is:

$$\text{Expected Time-to-Completion} = T_E = \frac{w}{Ca_c} + \frac{d_i + d_o}{Ba_b}$$

As noted earlier the client specifies the CPU speed,  $C_v$ , bandwidth,  $B_v$ , and reliability,  $r$ , of a virtual node. If reliability is not specified, the system of the present sets a default value for the baseline reliability  $R$ . The client may select a greater, but not lesser, reliability.

The system then computes the implied "requested" time-to-completion for the task as:

$$\text{Requested Time-to-Completion} = T_R = \frac{w}{C_v} + \frac{d_i + d_o}{B_v}$$

The subset of provider computers that satisfy the following criteria are considered "candidate provider computers":

- The provider computer is not working on a paying job.
- When the client specifies  $w$ ,  $d_i$ , and  $d_o$ , the provider computer has (approximately)  $T_E = T_R$ .
- When the client does not specify  $w$ ,  $d_i$ , and  $d_o$ , the provider computer has (approximately)  $a_c C = C_v$  and  $a_b B = B_v$ .

Taking provider computers' CPU's and CPU availabilities as fixed, candidate provider computers have expected bandwidths (bandwidth multiplied by bandwidth availability) that satisfy the following equation:

$$Ba_b = Ca_c \frac{(d_i + d_o)C_v B_v}{wB_v(Ca_c - C_v) + (d_i + d_o)Ca_c C_v}$$

For each candidate provider computer, the system calculates the weighted average availability,  $\pi$ , as:

1 computer's power rating. Note that, effectively, the provider is being paid for time, not work.  
 2 However, because the provider computer's power rating is a function of the provider  
 3 computer's availability, the rate of payment the provider receives is a function of the average  
 4 work per unit time the provider computer completes over time.

5 Providers are paid a fixed amount per unit time. The system established a target annual  
 6 payment per MHz-hour,  $k$ , and a baseline reliability,  $R$ . A provider computer with a weighted  
 7 average availability of  $\pi$  and CPU of  $C$  is paid  $\frac{kC \ln(1-\pi)}{\ln(1-R)}$  per hour. This assumes that the  
 8 provider computer is working over the whole interval from launch to conclusion of the task.

#### 9 Provider Payments and Task Scheduling (Full Version)

10

11 As part of establishing a given provider computer as a viable source of computing  
 12 power for the system, for each provider computer, the system constructs a CPU signature and  
 13 a bandwidth signature. These signatures show the probabilities of the provider computer's  
 14 CPU and bandwidth being available for a fixed time interval over time. For example:

15

Beginning of Time Interval	Pr(CPU available)	Pr(Bandwidth available)
1	30%	0%
2	40%	10%
3	50%	20%
4	40%	30%

16

17 The number of intervals should be enough to cover a reasonable circadian "cycle"  
 18 (e.g. 1 week). The signature should be updated periodically.

19 For each provider computer, the system measures CPU speed by sending to the  
 20 provider computer certain benchmark tasks that have known response times. The bandwidth  
 21 is also measured. As noted earlier, these measures are periodically updated.

1 described, multiple instances of each such module may, in fact, be present in the architecture  
2 in order to allow for scalability to larger operations, to increase reliability so that no single  
3 point of failure exists in the system, or to have multiple instances of the entire server located  
4 at different physical locations to allow for redundancy and insurance against failure.

5 The server comprises a registration servlet 50, which accepts registration requests  
6 over the network 28 from clients and providers who wish to participate on the system.  
7 Registration information is sent by the registration servlet to the registration manager 52 for  
8 processing.

9 The registration manager 52 receives appropriate information from the registration  
10 servlet 50 and assigns the client/provider computer to a database. The registration manager  
11 then generates a certificate and returns the certificate in the response to the client/provider  
12 computer. The provider computer certificate includes the provider public key, the provider-id  
13 (which is a hash of the provider public key), and the provider database ID. The client  
14 certificate comprises the client/user public key, which then acts as the client/user ID, the  
15 client ID, and the provider computer database ID.

16 The provider servlet 56 accepts various messages from providers as will be discussed  
17 more fully below, and routes them to the appropriate provider manager 58. Further, the  
18 provider servlet 56 returns responses from the provider manager to the various provider over  
19 the network 28. The web server of the present invention encrypts and decrypts requests,  
20 authenticates requests, routes them to the appropriate servlet, encrypts and returns responses.  
21 The provider servlet accepts GET and POST messages from providers. (These messages will  
22 be more fully explained below).

23 Provider manager 58 manages the interaction of the system of the present invention  
24 with the various providers. The provider manager receives and routes requests from other  
25 internal server components such as the client manager 68, the task scheduler 64, and the  
26 database server 54. All administrative information and processing that relates to providers is  
27 managed by the provider manager 58.

28 The element servlet 60 is used as the point where data elements and code are  
29 exchanged between clients and providers over the network. Elements are executables, or

1 in the database 54. Various algorithms present in the task scheduler ensure that appropriate  
2 tasks are assigned to the appropriate providers given priority of tasking and other prioritizing  
3 data.

4 Database server 54 stores information on clients and providers. To accomplish this  
5 task, the database server 54 comprises an account database relating to providers and revenue  
6 accumulated by providers, a client database relating to the identity and contact information for  
7 clients, job and task information, and a provider database relating to providers and  
8 identification and contact information for providers. Such information, while shown as being  
9 located on a single database server 54 can also be located on multiple servers to so that the  
10 system can be scaled up to handle many more users as well as ensure privacy and security of  
11 the information.

12 The web manager 70 intermediates access between the account section of the web site  
13 and the database. The web manager 70 supports the web based registration, update of  
14 registration and account information, and retrieval of account-specific and system-wide  
15 statistics.

16 Monitoring and control module 72 comprises a set of utilities that support the overall  
17 management of the system of the present invention. The monitoring and control functionality  
18 72 provides graphical user interfaces for system operators to monitor the health of the overall  
19 system. The monitoring and control functionality 72 polls all servers for health information in  
20 the form of name-value pairs that contains status information and related attributes to all of  
21 the servers. In this fashion, the monitoring and control 72 can keep track of the health of the  
22 entire system. The monitoring and control module can also respond to asynchronous  
23 notification of system events from servers.

24 Log server 74 contributes to the health of the overall system by collecting debug and  
25 error log data. Any errors that are encountered throughout the system whether they are in  
26 clients, providers, or any of the servers of the system are collected in the log server for  
27 subsequent analysis and debugging where necessary.

28 As noted above, while single instances of the various modules have been indicated in  
29 **Figure 5**, this architecture is not meant as a limitation. Additional instances of any given



1 accomplished in order to even the workflow and ensure the most efficient functioning of the  
2 system of the present invention.

3 In a similar fashion, web manager 70 can be augmented by additional instances of web  
4 managers with requests being allocated among the various web managers. The monitoring  
5 and control module 72 can also be represented by multiple instances of the monitoring and  
6 control function. There can also be multiple instances of the log server 74. However, there  
7 must be close coordination among the various log servers so that any trends and errors can be  
8 detected and corrected.

9 Registration servlet 50 can also be augmented by additional components. Clients and  
10 providers are assigned across the various instances of registration servlets.

12 The task scheduler 64 maintains an in-memory model of the state of all of the  
13 providers (i.e. the provider profiles) and any tasks that it manages. Information to establish  
14 the task scheduler is initially retrieved from database 54. It is thereafter updated as messages  
15 are received from the provider manager regarding status of efforts of the various providers,  
16 and from the client manger, which forwards relevant messages from clients to the task  
17 scheduler.

18 The task scheduler also requests additional information about providers by cuing a get  
19 task message (for retrieving task status) or a get cache contents for obtaining the cache  
20 contents of the server for the provider. The task scheduler comprises various scheduling  
21 algorithms. This task scheduler makes initial assignments of tasks to providers and migrates  
22 tasks from one provider to another if necessary. When the task scheduler schedules a new  
23 task, it records the assignment of the task and the provider to which it is assigned in the task  
24 database 54. It further queues a task message to the provider. When migrating a task, the task  
25 scheduler removes the task from the current provider by cuing a remove task message,  
26 assigning the partially completed task to a new provider via a task message. All such  
27 operations are reported in the task database.

28 The element servlet is the entity that stores data elements that are to be assigned to  
29 providers for processing. Clients upload their data elements to the element servlet 60. The  
30 element servlet authenticates the client as one that is permitted to store elements in the file

1 only access the server if they have a legitimate certificate. In fact, the server uses a hash of the  
2 client's public key as a client identifier.

3 The protocol is reliable because it incorporates sequence numbers, acknowledgement,  
4 and retransmission. Combined, these protect against dropped messages, out-of-order  
5 processing of messages, and the processing of duplicate messages.

6 It is worth noting, the low level protocol provides the abstraction of a bi-directional  
7 communication channel between components even though it is implemented on top of HTTP  
8 that has a request-response model. A bi-directional model is simulated by having the client  
9 periodically "poll" the server for messages. This poll is done on an exponential decay to limit  
10 the amount of traffic generated by polling. If the CTS does not respond to a particular GET  
11 message, perhaps because of heavy CTS load, subsequent GET messages are sent ever more  
12 infrequently, according to an exponential decay, until some limit periodicity is achieved. In  
13 this manner, heavy CTS can be automatically mitigated through less frequent CE GET  
14 requests.

15 The maximum polling rate is present each time the client and server exchange a high  
16 level message. The two basic operations in the low level protocol are the "GET" and the  
17 "POST", both of which are initiated by the client. The POST operation sends a new message  
18 (or set of messages) to the server. The GET is a poll operation to get any message (or set of  
19 messages) that the server may have enqueued for the client. Since the low-level protocol is  
20 based on HTTP it takes advantage of that protocol's features for multipart messages and  
21 different message encodings (e.g., for compression). Further, multipart messages can be sent  
22 in order to maximize communication efficiency. All messages are also subject to data  
23 compression to limit the volume of data being transmitted and to conserve bandwidth.

24 The low-level protocol also supports the concept of session. Sessions are implemented  
25 using cookies and provide for the resetting of a communication channel in the event that  
26 synchronization is lost.

27 The following is a listing of high level message types used by the present invention in  
28 its communications protocol. This list is illustrative in nature. Other message types can be  
29 added to the present invention as the need arises. As noted above, a series of messages are

- 1                   Request that a DataElement message be sent describing a particular data  
2                   element
- 3   RemoveExecutableElement:  
4                   Request the removal of the named executable element.
- 5   RemoveDataElement:  
6                   Request the removal of the named data element.
- 7
- 8   Attach:  
9                   Attach to a particular task or job. The attach message is used to subscribe to  
10                  receive status updates on the identified task or job.
- 11   GetTask:  
12                  Request that the identified task be checkpointed. The checkpoint is returned to  
13                  the client in a Task message.
- 14   GetTaskStatus:  
15                  Requests the status of a particular task. The server will send a TaskStatus  
16                  message in response.
- 17   RemoveTask:  
18                  Requests the removal of the named task.
- 19   ExecutableElement:  
20                  Register an executable element with the server. The actual executable element  
21                  data must be downloaded over the data channel.
- 22   DataElement:  
23                  Register a data element with the server. The actual data element data must be  
24                  downloaded over the data channel.
- 25   CloseSession:  
26                  Close current session.
- 27   Error:  
28                  Report an error condition.
- 29   ExternalDelivery:

- 1 Alert the recipient that a message is waiting at a specified URL.
- 2
- 3 Provider-to-Server messages
- 4 NodeStatus:
- 5 Reports a change in the status of the provider node.
- 6 NodeProfile:
- 7 Reports profile information for node (e.g., machine type, available disk,
- 8 available memory, connection speed, operating system type).
- 9 GetConfig:
- 10 Request that the server send a Config message to set the providers server-
- 11 controlled configuration settings.
- 12 CacheContents:
- 13 Reports the contents of the providers element cache. This information is useful
- 14 in task scheduling.
- 15 Contents:
- 16 Reports the contents of the task queue.
- 17 Task Status:
- 18 Reports the status of a task. There can be both intermediate and final task
- 19 statuses. Intermediate task statuses may contain partial results. Final task
- 20 statuses contain the results of the task. Task statuses are also used to report
- 21 errors in the execution of a task.
- 22 Task: The provider sends the task message out in response to a GetTask message
- 23 from the provider. In this case, the Task message is a checkpoint of the state of
- 24 the identified task.
- 25 Error:
- 26 Report an error condition.
- 27 ExternalDelivery:
- 28 Alert the recipient that a message is waiting at a specified URL.
- 29 Heartbeat:

1           The present invention can also be employed in other manners, such as a method of  
2 marketing computers by offering incentives to computer customers that agree to operate a  
3 compute engine (CE) on the computers and having the CE utilize idle computational  
4 processing power on the computers. Incentives can include, but are not limited to free  
5 computer use, free ISP service, discounted computer sales price, discount computer lease  
6 price, a sales rebate, periodic rebates, and usage fees. The CE can also be utilized to deliver  
7 "pushed" content, such as advertising, to these computer customers via a display window of  
8 said computer's graphic interface or via said computer's sound output.

9           A system and method for the establishment and utilization of networked idle  
10 computational processing power has been illustrated. It will be appreciated by those skilled in  
11 the art that other variations in, for example, the calculation of power, and the methods for  
12 compensating providers can be accomplished without departing from the scope of the  
13 invention as disclosed.

14

- 1 8. The system for using excess computational power of claim 5 wherein the run time  
2 interface further comprises instructions for attaching the client computer of a job in progress  
3 and for checking on the status of a job while the job is in progress.
- 4 9. The system for using excess computational power of claim 1 wherein the CTS further  
5 comprises instructions for charging the at least one client for computational power used  
6 during the execution of the programs.
- 7 10. The system for using excess computational power of claim 9 wherein the charging of  
8 the client is made when the task is launched.
- 9 11. The system for using excess computational power of claim 10 wherein a task is  
10 launched when it is scheduled by the CTS to be run by at least one provider computer..
- 11 12. The system for using excess computational power of claim 1 wherein the CTS further  
12 comprises instructions for monitoring a virtual node specified by the client for failure.
- 13 13. The system for using excess computational power of claim 9 wherein charging of  
14 clients is based upon the sum of the power of the provider computers running the client tasks.
- 15 14. The system for using excess computational power of claim 9 wherein the CTS further  
16 comprises instructions for making payments to providers based upon the available  
17 computational power the provider's computer.
- 18 15. The system for using excess computational power of claim 14 wherein the available  
19 computational power is a function of CPU availability, CPU speed and bandwidth  
20 availability.
- 21 16. The system for using excess computational power of claim 15 wherein CPU speed is  
22 measured by running benchmark tests of the CPU periodically.
- 23 17. The system for using excess computational power of claim 1 wherein each of the  
24 plurality of provider computer further comprises a general purpose compute engine (CE) for  
25 executing programs sent from the central server.
- 26 18. The system for using excess computational power of claim 17 wherein the CE further  
27 comprises instructions for allowing an owner of the provider computer to set the availability  
28 parameters of the provider computer.

- 1 25. The system for using excess computational power of claim 24 wherein the client is  
 2 also charged an hourly premium for requesting reliability greater than a set baseline reliability  
 3 given by the formula:

$$4 \quad \text{Price per hour for premium reliability} = kCN \left( \frac{\ln(1-r)}{\ln(1-R)} - 1 \right)$$

- 5  
 6 where k is a baseline price per MHz-hour, C is the CPU speed requested of a virtual node  
 7 requested by the client, N is the number of virtual nodes running, r is the requested reliability,  
 8 and R is the baseline reliability.

- 9 26. The system for using excess computational power of claim 1 wherein the owners of  
 10 the provider computers are paid a fixed amount per unit work for providing excess  
 11 computational power according to the following formula:

$$12 \quad \frac{T \ln(1-\pi)}{\ln(1-r)} \text{ per GHz-hour, where}$$

- 13 T is a rate,  $\pi$  is an average availability, and r is a target reliability.

- 14 27. A method for using computer excess computational power comprising:  
 15 adapting at least one client computer to create programs for execution;  
 16 connecting the client computer to a network;  
 17 connecting a central task server (CTS) to the network, wherein said CTS is adapted to  
 18 receive the programs from the client;  
 19 connecting a plurality of provider computers for providing excess computational  
 20 power to the network and adapting said provider computers to receive the programs from the  
 21 central server, the plurality of provider computers providing computational power upon  
 22 demand.

- 23 28. The method for using excess computational power of claim 27 further comprising  
 24 supplying the at least one client computer with a client application programming interface  
 25 (API) for allowing the client computer to specify the parameters under which the programs  
 26 are to be run.

- 1 40. The method for using excess computational power of claim 35, further comprising  
2 providing the CTS with instructions for making payments to providers based upon the  
3 available computational power the provider's computer.
- 4 41. The method for using excess computational power of claim 40 wherein the available  
5 computational power is determined as a function of CPU availability, CPU speed and  
6 bandwidth availability.
- 7 42. The method for using excess computational power of claim 41 wherein CPU speed is  
8 measured by running benchmark tests of the CPU periodically.
- 9 43. The method for using excess computational power of claim 27, further comprising  
10 providing each of the plurality of provider computer with a general purpose compute engine  
11 (CE) for executing programs sent from the central server.
- 12 44. The method for using excess computational power of claim 43, further comprising  
13 providing the CE with instructions for allowing an owner of the provider computer to set the  
14 availability parameters of the provider computer.
- 15 45. The method for using excess computational power of claim 44 wherein the availability  
16 parameters comprise the amount of RAM available for the programs, the amount of disk  
17 space available for the programs and the times during the day, if any, that the provider  
18 computer is available to execute the programs.
- 19 46. The method for using excess computational power of claim 44, further comprising  
20 providing the CE with instructions for opening a window on the provider computer when the  
21 programs are being run.
- 22 47. The method for using excess computational power of claim 44, further comprising  
23 providing the CE with instructions for displaying content from the central server in the  
24 window.
- 25 48. The method for using excess computational power of claim 44, further comprising  
26 providing the CE with instructions for executing the programs only when excess  
27 computational power is detected on the provider computer.



1         $T$  is a rate,  $\pi$  is an average availability, and  $r$  is a target reliability.

2        53.     A method for the establishment and utilization of networked idle computational  
3        processing power comprising:

4                loading a client API on to client computers to allow client jobs to be broken into  
5        independent tasks.

6                loading compute engine (CE) software onto provider computers to execute tasks using  
7        idle computational processing power of said provider computers and using a centralized task  
8        server (CTS) that is at least intermittently networked to said client computers at any time  
9        selected by said clients and intermittently connected to provider computers to exchange tasks  
10       and results.

11       54.     The method of claim 53, wherein said CE runs unobtrusively and securely on the  
12       provider computer.

13       55.     The method of claim 53, wherein a task scheduler in said CTS matches  
14       computational power requested by clients to provider resources.

15       56.     The method of claim 53, wherein CE operates the provider computer to  
16       periodically contact the CTS to retrieve tasks and control messages and/or send results.

17       57.     The method of claim 53, wherein the client API periodically contacts the CTS to  
18       retrieve results.

19       58.     The method of claim 53, wherein the CE displays a graphic window when running  
20       tasks to display processing and status information.

21       59.     The method of claim 58, further comprising display of advertising on other  
22       targeted messages by the CE graphic window.

23       60.     The method of claim 58, wherein the window can be made to disappear based on  
24       configurable mouse coordinates.

25       61.     A system for the established and utilization of networked idle computational  
26       processing power comprising:

27                at least one client computer with a client API that allows client jobs to be broken into  
28       independent tasks;

- 1 charging said clients a fee; and  
2 compensating said providers.
- 3 68. The method of creating and intermediating a market for idle computational  
4 processing power of claim 67 wherein API software on client computers provides  
5 tasks and compute engine (CE) software on said provider computers executes said  
6 tasks and provides said results.
- 7 69. The system of claim 61, further comprising a communication protocol for  
8 governing communication between the API, CTS, and CE, said protocol having  
9 separate message and data channels; wherein the message channel is a bi-  
10 directional channel between clients and the server that is used to pass high-level  
11 messages that facilitate distributed processing and the data channel is a channel  
12 used to pass large blocks of data elements and executable elements.
- 13 70. The system of claim 61, wherein the client API supports jobs of arbitrary size by  
14 being programmed to monitor only task results.
- 15 71. The system of claim 61, further comprising a runtime API which defines how  
16 tasks interact with their environment.
- 17 72. The system of claim 61, wherein the client API contains an interface for task  
18 checkpointing.
- 19 73. The system of claim 61, wherein the CE supports the queuing of tasks, task  
20 prioritization, and preemption of tasks based on priority.
- 21 74. The system of claim 61, wherein the CTS and CE together record computer CPU  
22 and communications line signatures.
- 23 75. The system of claim 74, wherein the signatures are historical usage patterns for use  
24 in task scheduling.
- 25 76. The system of claim 61, further comprising system-wide secure communications,  
26 client anonymity, and obfuscation of client code.
- 27 77. The system of claim 61, further comprising means for pricing jobs and scheduling  
28 tasks based on delivering clients a virtual cluster comprised of virtual nodes, said

1/8

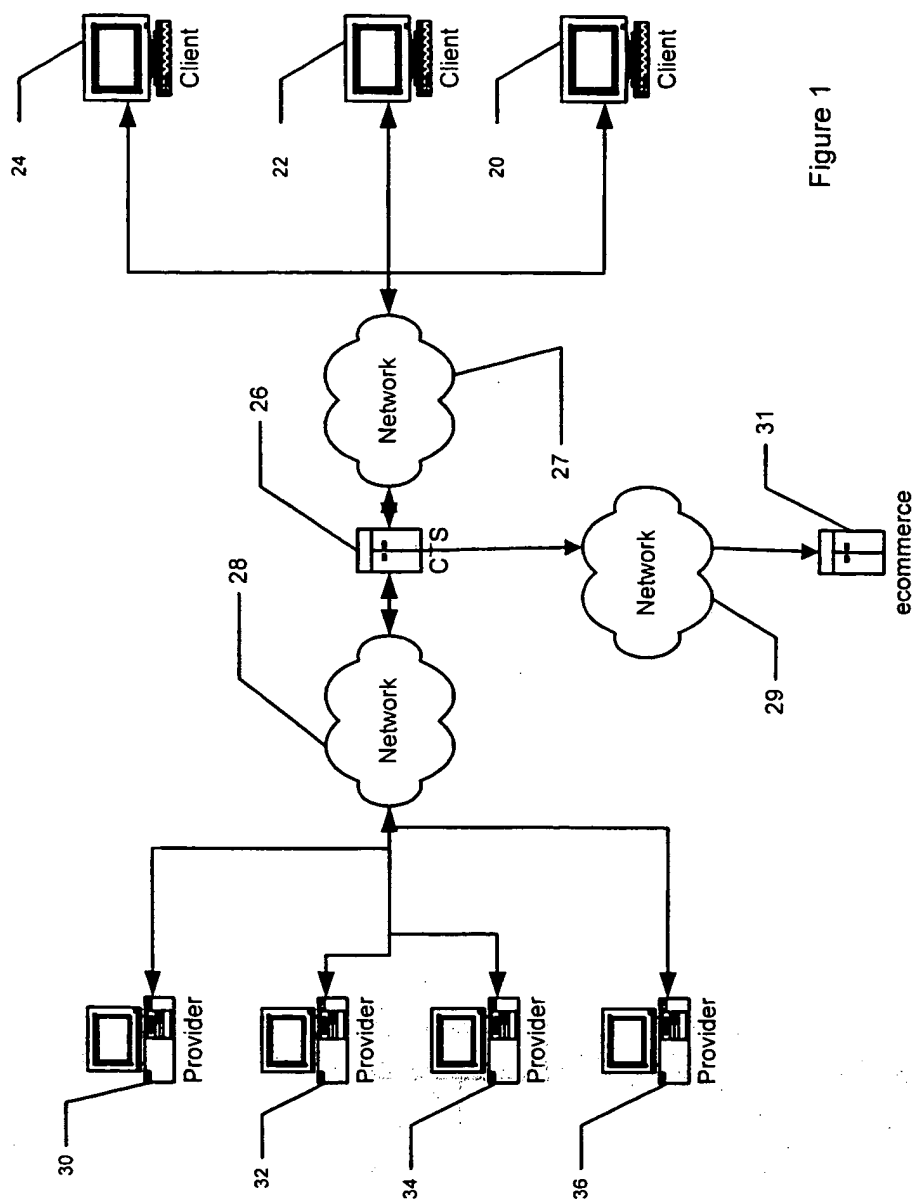
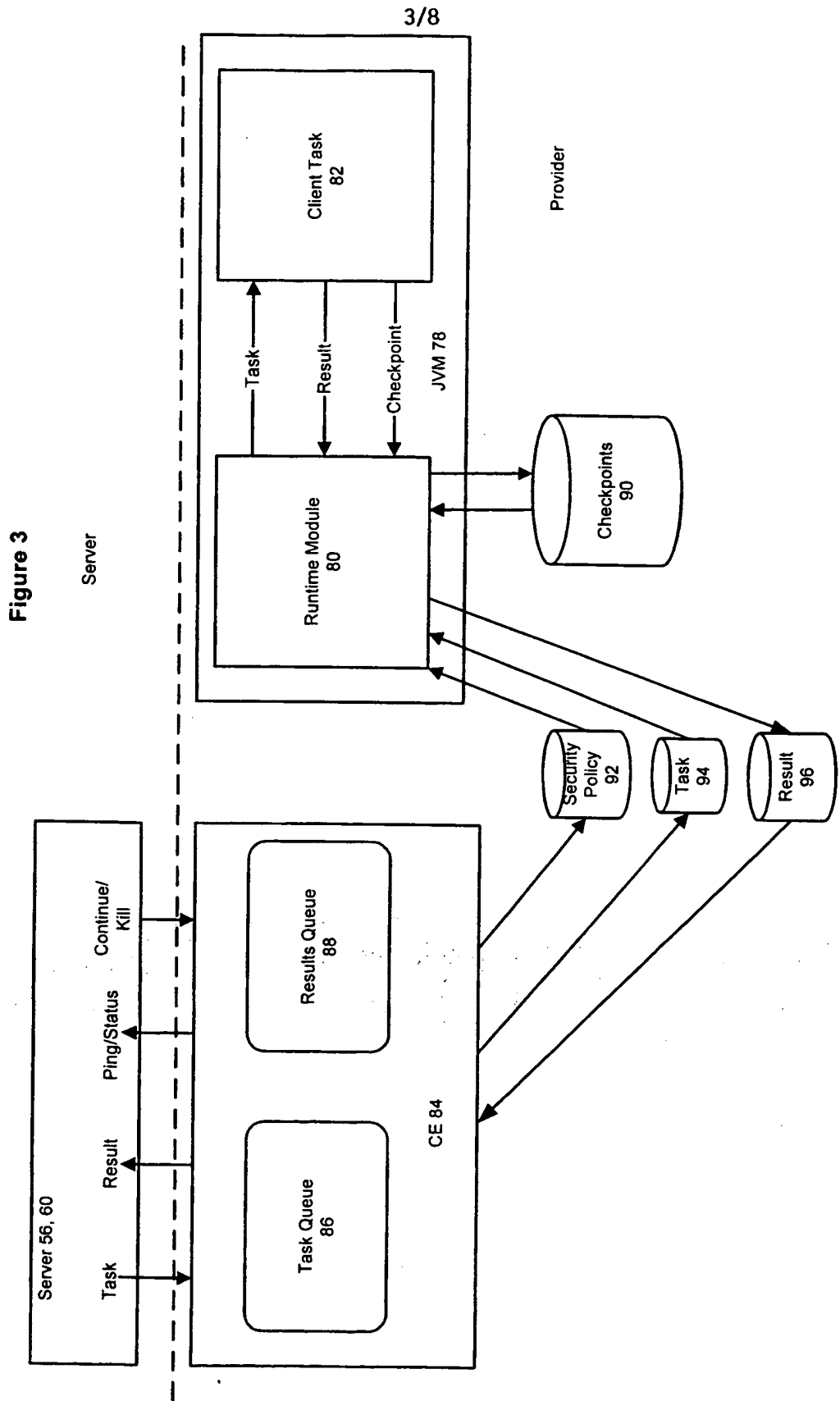


Figure 1



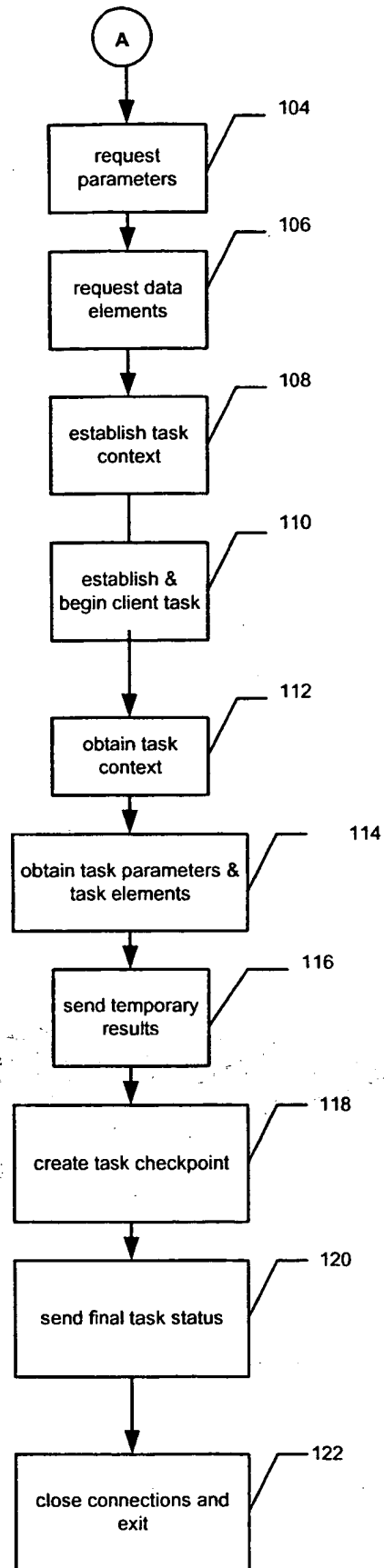
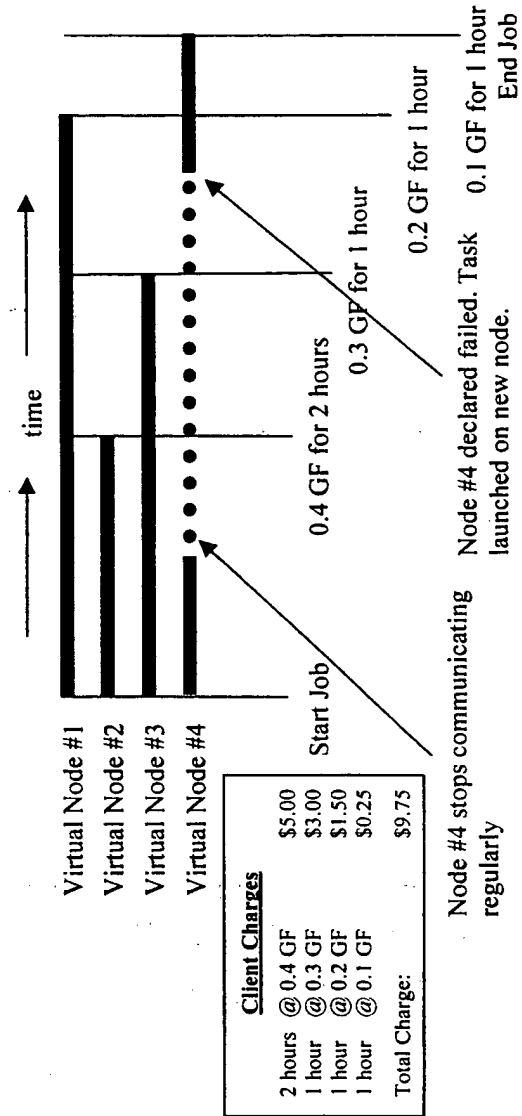


Figure 4B

7/8

Figure 6



(19) World Intellectual Property Organization  
International Bureau(43) International Publication Date  
1 March 2001 (01.03.2001)

PCT

(10) International Publication Number  
**WO 01/14961 A3**(51) International Patent Classification<sup>7</sup>: G06F 9/50

(21) International Application Number: PCT/US00/24336

(22) International Filing Date: 28 August 2000 (28.08.2000)

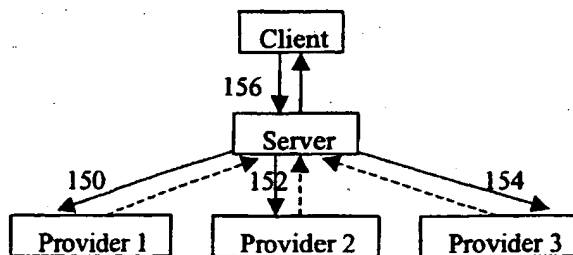
(25) Filing Language: English

(26) Publication Language: English

(30) Priority Data:  
60/150,766 26 August 1999 (26.08.1999) US  
60/210,334 12 June 2000 (12.06.2000) US(71) Applicant: PARABON COMPUTATION [US/US];  
Suite 1000, 3930 Walnut Street, Fairfax, VA 22030 (US).(72) Inventors: ARMENTROUT, Steven, L.; 11879 St.  
Trinians Ct., Reston, VA 20191 (US). O'CONNOR,  
James; 11625 Ayreshire Road, Oakton, VA 22124 (US).  
GANNON, James; 3817 Inverness Road, Fairfax, VA  
22030 (US). SLETTEN, Brian; Apt. 906, 900 N. Stuart  
Street, Arlington, VA 22203 (US). CIER, Sean; Apt.  
12, 11800 Federalist Way, Fairfax, VA 22030 (US).  
CARLSON, Sarah; 13236 Memory Lane, Fairfax, VA  
22033 (US). DAVIS, Jonathan; Apt. 312, 202 Columbia  
Road NW, Washington, DC 20009 (US). DUPERTUIS,  
Greg; 7888 Leroux Lane, Manassas, VA 20112 (US).  
MCCLOUGHLIN, Scott; 2827 28th Street, #14, Wash-  
ington, DC 20001 (US). DAVIES, Antony; 11733 Stuart  
Mill Road, Oakton, VA 22124 (US).(74) Agents: ROBERTS, Jon, L. et al.; Roberts Abokhair &  
Mardula, LLC, Suite 1000, 11800 Sunrise Valley Drive,  
Reston, VA 20191 (US).(81) Designated States (*national*): AE, AG, AL, AM, AT, AU,  
AZ, BA, BB, BG, BR, BY, BZ, CA, CH, CN, CR, CU, CZ,  
DE, DK, DM, DZ, EE, ES, FI, GB, GD, GE, GH, GM, HR,  
HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR,  
LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, MZ,  
NO, NZ, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM,  
TR, TT, TZ, UA, UG, UZ, VN, YU, ZA, ZW.(84) Designated States (*regional*): ARIPO patent (GH, GM,  
KE, LS, MW, MZ, SD, SL, SZ, TZ, UG, ZW), Eurasian  
patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European  
patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE,  
IT, LU, MC, NL, PT, SE), OAPI patent (BF, BJ, CF, CG,  
CI, CM, GA, GN, GW, ML, MR, NE, SN, TD, TG).

Published:

— with international search report

(88) Date of publication of the international search report:  
2 August 2001For two-letter codes and other abbreviations, refer to the "Guid-  
ance Notes on Codes and Abbreviations" appearing at the begin-  
ning of each regular issue of the PCT Gazette.(54) Title: SYSTEM AND METHOD FOR THE ESTABLISHMENT AND UTILIZATION OF NETWORKED IDLE COMPU-  
TATIONAL PROCESSING POWER(57) Abstract: A distributed computing platform  
using the idle computational processing power of  
a plurality of provider computers is disclosed. At  
least one networked server collects tasks from client  
computers, schedules and distributes the tasks to  
networked provider computers, and collects and  
returns results to client computers. A client API forms  
tasks and collects results. A compute engine operates  
on the provider computers to communicate with the  
server and execute tasks using idle computational  
power.

WO 01/14961 A3

## C.(Continuation) DOCUMENTS CONSIDERED TO BE RELEVANT

Category *	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	CHAO-JU HOU ET AL: "IMPLEMENTATION OF DECENTRALIZED LOAD SHARING IN NETWORKED WORKSTATIONS USING THE CONDOR PACKAGE" JOURNAL OF PARALLEL AND DISTRIBUTED COMPUTING,US,ACADEMIC PRESS, DULUTH, MN, vol. 40, no. 2, 1 February 1997 (1997-02-01), pages 173-184, XP000682281 ISSN: 0743-7315 page 174, left-hand column, line 35 -page 175, right-hand column, last line ---	1-3,17, 27-29, 43,63, 79,80
A	WALDSPURGER C A ET AL: "Spawn: a distributed computational economy" IEEE TRANSACTIONS ON SOFTWARE ENGINEERING,US,IEEE INC. NEW YORK, vol. 18, no. 2, February 1992 (1992-02), pages 103-117, XP002124500 ISSN: 0098-5589 the whole document ---	1-81
A	MUTKA M W: "ESTIMATING CAPACITY FOR SHARING IN A PRIVATELY OWNED WORKSTATION ENVIRONMENT" IEEE TRANSACTIONS ON SOFTWARE ENGINEERING,US,IEEE INC. NEW YORK, vol. 18, no. 4, 1 April 1992 (1992-04-01), pages 319-328, XP000328663 ISSN: 0098-5589 the whole document ---	67,74
A	TANDIARY F ET AL: "Batrun: utilizing idle workstations for large scale computing" IEEE PARALLEL AND DISTRIBUTED TECHNOLOGY: SYSTEMS AND APPLICATIONS,XX,XX, vol. 4, no. 2, 21 June 1996 (1996-06-21), pages 41-48, XP002124512 ISSN: 1063-6552 the whole document -----	1,17,27, 43,53, 54,56,61